

lecture 3

- optimal h and bias-Var for K-D
- KDE

L^{III} Kernel, m
density est

- OH MMP
- Tue 4/4 at 4pm
CSSS Conference
Room
PDL level LL

JF Office h
Mon 1:30 - 2:30

CSSS Conf.
Room

Lecture II – Nearest Neighbor and Kernel predictors

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

STAT/BIOST 527
Spring 2023

① Nearest-Neighbor predictors ✓

② Kernel predictors ✓

③ An elementary analysis of Kernel Regression ↙

④ Bias, Variance and h for $x \in \mathbb{R}$ ↙

h^* theoretical

in practice: CV next lecture

Regression example by N-W

A special case in wide use is the Nadaraya-Watson regressor

$$\hat{y}(x) = \frac{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right) y^i}{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right)} = \sum w_i b_h(x, x^i) \quad (4)$$

In this regressor, $f(x)$ is always a convex combination of the y^i 's, and the weights are proportional to $b_h(x, x^i)$.

The Nadaraya-Watson regressor is biased if the density of P_x varies around x .

Assumed $y(x) = f(x) + \varepsilon$ true !! noise

$$\hat{y}(x) - y(x) = \text{error}$$

depends on ε

\uparrow
depends on $x^{1:n}, y^{1:n}$ data

Kernel regression by Nadaraya-Watson

$$\hat{y}(x) = \frac{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right) y^i}{\sum_{i=1}^n b\left(\frac{\|x-x^i\|}{h}\right)} \quad (6)$$

Let $w_i = \frac{b\left(\frac{\|x-x^i\|}{h}\right)}{\sum_{i'=1}^n b\left(\frac{\|x-x^{i'}\|}{h}\right)}$.

Assumptions

- A0 For simplicity, in this analysis we assume $x \in \mathbb{R}$.
 A1 There is a true smooth¹ function $f(x)$ so that

$$y = f(x) + \varepsilon, \quad (7)$$

where ε is sampled independently for each x from a distribution P_ε , with $E_{P_\varepsilon}[\varepsilon] = 0$, $Var_{P_\varepsilon}(\varepsilon) = \sigma^2$.

- A2 The kernel $b(z)$ is smooth, $\int_{\mathbb{R}} b(z) dz = 1$, $\int_{\mathbb{R}} z b(z) dz = 0$, and we denote $\sigma_b^2 = \int_{\mathbb{R}} z^2 b(z) dz$, $\gamma_b^2 = \int_{\mathbb{R}} b^2(z) dz$.

In this first analysis, we consider that the values $x, x^{1:N}$ are fixed; hence, the randomness is only in $\varepsilon^{1:N}$.

¹with continuous derivatives up to order 2

Expectation of $\hat{y}(x)$ – a simple analysis

Expanding f in Taylor series around x we obtain

$$f(x^i) = f(x) + f'(x)(x^i - x) + \frac{f''(x)}{2}(x^i - x)^2 + o((x^i - x)^2) \quad (8)$$

We also have

$$y^i = f(x^i) + \varepsilon^i. \quad (9)$$

We now write the expectation of $\hat{y}(x)$ from (6), replacing in it y^i and $f(x^i)$ as above. What we would like to happen is that this expectation equals $f(x)$. Let us see if this is the case.

$$E_{P_\varepsilon^n} [\hat{y}(x)] = E_{P_\varepsilon^n} \left[\sum_{i=1}^n w_i y^i \right] = E_{P_\varepsilon^n} \left[\sum_{i=1}^n w_i (f(x^i) + \varepsilon^i) \right] \quad (10)$$

$$= \sum_{i=1}^n w_i f(x) + \sum_{i=1}^n w_i f'(x)(x^i - x) + \sum_{i=1}^n w_i \frac{f''(x)}{2}(x^i - x)^2 + E_{P_\varepsilon^n} \left[\underbrace{\sum_{i=1}^n w_i \varepsilon^i}_{=0} \right] \quad (11)$$

Vary

$$= f(x) + f'(x) \underbrace{\sum_{i=1}^n w_i (x^i - x)}_P \quad \text{bias} \quad (12)$$

Px

$$+ \frac{f''(x)}{2} \underbrace{\sum_{i=1}^n w_i (x^i - x)^2}_\text{bias}$$

In the above, the expressions in red depend of f and x , those in blue depend on $x^{1:n}$.

given $\mathcal{D} = (x^{1:n}, y^{1:n})$, x

$$E[\hat{y}(x) \mid \mathcal{D}, f, h, \text{noise}(0, \sigma^2)] = f(x) + \underbrace{\text{bias}_h(f, P_x)}_{\text{prediction}} + \text{bias}_w(f'')$$

prediction

$$\text{bias } E[\hat{y}(x) - f(x)] =$$

$$\text{Var } \hat{y}(x) \mid \mathcal{D}, f, h, \text{noise}$$

- MSE $\underline{\hat{y}(x)} = \underline{\text{bias}^2(\hat{y}(x))} + \underline{\text{Var } \hat{y}(x)}$

Mean
squared Error

average over \mathcal{D} : $x^{1:n} \sim P_x$

$$y = f(x) + \varepsilon$$

- MSE $[\hat{y}(x) \mid P_x, f, h, w, \text{noise } \sigma^2]$

- MISE = $E_{P_x} [\text{MSE } \hat{y}(x)] = \int_{\mathbb{R}} \text{MSE}[\hat{y}(x)] P_x dx$

Integrated

Bias, Variance and h for $x \in \mathbb{R}$

2

The **bias** of \hat{y} at x is defined as $E_{P_X^n} E_{P_\varepsilon^n} [\hat{y}(x) - f(x)]$.

bias

$$E_{P_X^n} E_{P_\varepsilon^n} [\hat{y}(x) - f(x)] = h^2 \sigma_b^2 \left(\frac{f'(x)p'_X(x)}{p_X(x)} + \frac{f''(x)}{2} \right) + o(h^2) \quad (13)$$

The **variance** \hat{y} at x is defined as $Var_{P_X^n P_\varepsilon^n}(\hat{y}(x))$.

$$\text{var} = \int b^2(\mathbf{x}) d\mathbf{x}$$

$$Var_{P_X^n P_\varepsilon^n}(\hat{y}(x)) = \frac{\gamma_b^2}{nh} \sigma^2 + o\left(\frac{1}{nh}\right). \quad (14)$$

The **MSE (Mean Squared Error)** is defined as $E_{P_X^n} E_{P_\varepsilon^n} [(\hat{y}(x) - f(x))^2]$, which equals

$$MSE(x) = \text{bias}^2 + \text{variance} = h^4 \sigma_b^4 \left(\frac{f'(x)p'_X(x)}{p_X(x)} + \frac{f''(x)}{2} \right)^2 + \frac{\gamma_b^2}{nh} \sigma^2 + \dots \quad (15)$$

$$MSE = \int MSE p_X dx =$$

²After

Optimal selection of h

$$\text{MISE}_{(h)} = h^4 + \frac{1}{nh}$$

If the MSE is integrated over \mathbb{R} we obtain the $\text{MISE} = \int_{\mathbb{R}} \text{MSE}(x) p_X(x) dx$.

The kernel width h can be chosen to minimize the MISE, for fixed f, p_X and b .
We set to 0 the partial derivative

find
min MISE

$$\frac{\partial \text{MISE}}{\partial h} = h^3 \left(\text{[blue square]} \right) - \frac{(\text{[orange square]})}{nh^2} = 0. \quad (16)$$

It follows that $h^5 \propto \frac{1}{n}$, or

$$h \propto \frac{1}{n^{1/5}}.$$

$$n^{-1} = 100,000n$$

$$h' = \frac{h}{10} \quad (17)$$

In d dimensions, the optimal h depends on the sample size n as

$$h \propto \frac{1}{n^{1/(d+4)}} \Rightarrow n^{-\left(1 - \frac{1}{d+4}\right)} = n^{-\frac{4}{5}} \quad (18)$$

The MSE with optimal h decreases as $\sim \frac{1}{n} n^{1/(d+4)}$

Compare this with the MSE of the mean of a distribution, which decreases $\sim \frac{1}{n}$

$$\text{Ex: } h \sim \frac{1}{n^\alpha} \quad \alpha > 0$$

MISE rate?

$$\bar{x} = \frac{1}{n} \sum x^i \quad \text{estimates } \bar{E}[x] \quad p_X$$

$$\text{MISE} \sim \frac{1}{n^\alpha}$$

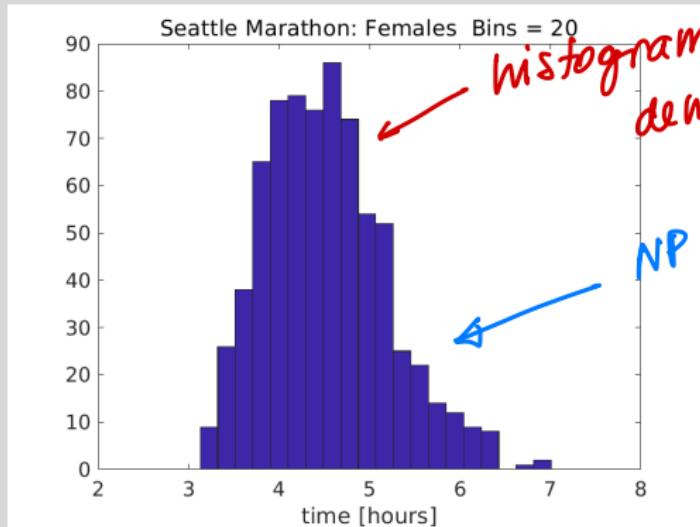
Lecture III – Kernel and k-NN density estimation

Marina Meilă
mmp@stat.washington.edu

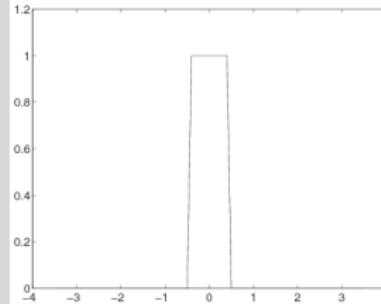
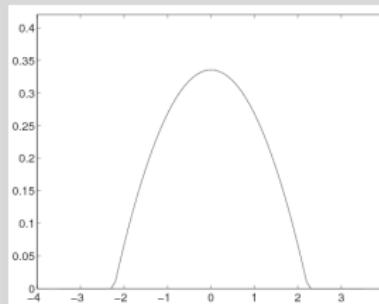
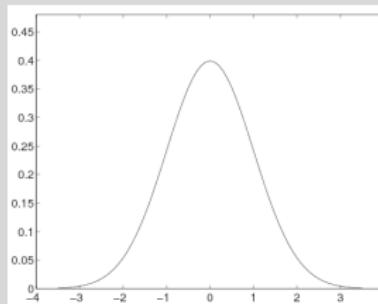
Department of Statistics
University of Washington

STAT/BIOST 527
Spring 2023

Arbitrary shaped densities



Kernels



$$E[b(z)] = 0$$

$$\int_{\mathbb{R}} b(z) dz = 1 \quad . \quad \text{a density!}$$

$$b \geq 0$$

also, usually $\int_{\mathbb{R}} z^2 b(z) dz = 1$

different b's can use same h

Kernel density estimators

$$\hat{P}_h(x) = \frac{1}{n} \sum_{i=1}^n b_h(x, x_i)$$

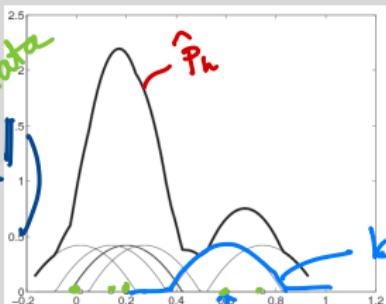
$$= \frac{1}{nh} \sum_{i=1}^n b\left(\frac{\|x - x_i\|}{h}\right)$$

$$\mathcal{F}_{h,b} = \{ \hat{P}_h \}$$

h = a smoothness parameter

$h \uparrow \hat{P}_h$ smoother

Var
bias ↗



$$\frac{1}{h} b\left(\frac{\|x - x_i\|}{h}\right)$$

$\underbrace{\qquad\qquad\qquad}_{z}$

$$\int_h b(z) dz = 1$$

for all $h >$

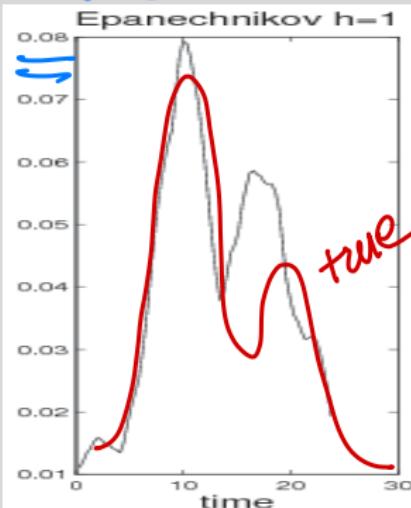
change of variable

$$\int_{-\infty}^{\infty} b(z) dz = 1$$

$$z = \frac{x - x_i}{h}$$

Choosing h

0.08



0.12

